

(51) Int.Cl.⁴
H 0 4 L 12/56

識別記号

F I
H 0 4 L 11/20

1 0 2 C

審査請求 未請求 請求項の数20 O L (全 11 頁)

(21) 出願番号 特願平10-27551
(22) 出願日 平成10年(1998) 2月9日
(31) 優先権主張番号 60/037843
(32) 優先日 1997年2月7日
(33) 優先権主張国 米国 (U S)
(31) 優先権主張番号 08/961122
(32) 優先日 1997年10月30日
(33) 優先権主張国 米国 (U S)

(71) 出願人 598077259
ルーセント テクノロジーズ インコーポ
レイテッド
Lucent Technologies
Inc.
アメリカ合衆国 07974 ニュージャージ
ー、マレーヒル、マウンテン アベニュー
600-700
(72) 発明者 アビジット カマー チョーデューリー
アメリカ合衆国, 07920 ニュージャージ
ー、スコッチ プレインズ、ハンティン
トン ロード 77
(74) 代理人 弁理士 三俣 弘文

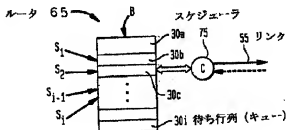
最終頁に続く

(54) 【発明の名称】 TCP接続の性能改善方法

(57) 【要約】

【課題】 フィードバック制御 TCP ネットワークの性能を改善する。

【解決手段】 フロー毎にキューイングする装置であり、複数の待ち行列に区分けされた所定サイズのバッファと、所定速度に従って各バッファからパケットを除去し、ネットワークを介して前記パケットを伝送するスケジューラと、受信パケットを受信することができ、待ち行列が使用可能である場合、パケットを待ち行列に入力することができるバッファ内の待ち行列の使用可能性を決定するためのコントロールデバイスとからなり、コントロールデバイスは更に、最長待ち行列ファーストスキームに従って待ち行列を選択し、待ち行列が使用可能でない場合、前記受信パケットの入力を可能にするために前記選択された待ち行列からパケットをドロップし、これにより、リンクによる公平性を高め、パケットスループットを高める。



【特許請求の範囲】

【請求項1】 (i)所定サイズのバッファを複数個の待ち行列に区分けするステップと、

ここで、各待ち行列は、情報パケットを受信し、一時的に格納するための占有率 b_i が割り当てられ、また、各バッファからパケットを除去し、そして、前記パケットをTCP接続を介して伝送するためのスケジューラにより処理される、

(ii)パケットが到着すると、前記パケットを受信するための待ち行列の使用可能性を決定し、前記待ち行列が使用可能である場合、前記パケットを待ち行列に入力するステップと、

(iii)前記待ち行列が使用可能でない場合、前記パケットの入力を可能にするため、待ち行列を選択し、そして、前記選択された待ち行列からパケットを解放するステップとからなり、これによりTCP接続の利用率を高めることを特徴とするTCP接続の性能改善方法。

【請求項2】 パケット到着時又は離脱時に、各待ち行列の現在の長さ q_i をトラッキングするステップを更に有する請求項1の方法。

【請求項3】 トラッキングステップは、パケットが前記待ち行列に入力されるたびに、前記待ち行列に付随するカウンタを増分するか、又は、パケットが前記待ち行列から解放されるときに、前記カウンタを減分するステップを有する請求項2の方法。

【請求項4】 前記スケジューラは所定の速度で前記各バッファからパケットを除去することができる請求項1の方法。

【請求項5】 前記バッファ待ち行列が満杯の場合、前記満杯待ち行列に宛てられたパケットが到着すると、満杯ではない1つ以上の待ち行列からバッファスペースを一時的に借りるステップを有し、前記満杯待ち行列の前記現在の長さはその名目的な占有率 b_i を超える請求項2の方法。

【請求項6】 待ち行列の使用可能性を決定する前記ステップは、所定サイズの前記バッファが満杯であるかどうかを決定するステップを有する請求項5の方法。

【請求項7】 前記選択された待ち行列から解放されるべき前記パケットは、前記待ち行列に内在する最も古いパケットである請求項2の方法。

【請求項8】 待ち行列を選択する前記ステップは、

(a)各待ち行列について現在の割当て b_i を確立するステップと、

(b)前記各待ち行列の現在の待ち行列長さの値 q_i を取得するステップと、

(c)現在の待ち行列長さの値 q_i と、各待ち行列の割当てられたバッファ占有率 b_i との間の差を計算するステップと、

(d)計算された差で最も大きな値を有する前記待ち行列を選択するステップとを有する請求項2の方法。

【請求項9】 待ち行列を選択する前記ステップは、

(a)各待ち行列について現在の割当て b_i を確立するステップと、

(b)現在の待ち行列長さの値 q_i が各待ち行列の割当てられたバッファ占有率 b_i を超える一連の待ち行列を計算するステップと、

(c)前記一連の待ち行列からランダムに待ち行列を選択するステップとを有する請求項2の方法。

【請求項10】 前記スケジューラは異なる所定速度に従って各バッファからパケットを除去し、待ち行列を選択する前記ステップは、

(a)そのそれぞれの待ち行列長さ q_i とその所定のサービス速度から各待ち行列に格納されたパケットにより経験されるキューイング遅延を計算するステップと、

(b)最も長いキューイング遅延を経験する待ち行列を選択するステップとを有する請求項2の方法。

【請求項11】 前記選択ステップは、

(a)1つ以上のビンを確立するステップと、ここで各ビンは所定の長さの1つ以上の待ち行列に結合される、

(b)最長長さを有する1つ以上の待ち行列に結合されたビンから待ち行列をランダムに選び取るステップを更に有する請求項2の方法。

【請求項12】 複数のソースからの情報パケットをTCP/IPネットワーク内の単一の通信リンクへ送信するルータであり、

(a)複数個の待ち行列に区分けされた所定サイズのバッファと、ここで、各待ち行列は、情報パケットを受信し、一時的に格納するための占有率 b_i が割り当てられる、

(b)各バッファからパケットを除去し、また、前記接続を介して前記パケットを伝送するためのスケジューラと、

(c)前記バッファの前記待ち行列が使用可能である場合、前記待ち行列に受信バッファを入力し、更に、前記バッファの前記待ち行列が使用可能でない場合、待ち行列を選択し、前記受信パケットの入力を可能にするため、前記スケジューラが前記選択された待ち行列からパケットを解放できるようにするための、前記バッファ内の待ち行列に使用可能性を決定するコントロール手段と、からなることを特徴とするルータ。

【請求項13】 パケットが入力されるたびに、又は前記待ち行列から解放されるたびに、前記待ち行列の現在の長さ q_i をトラッキングするための、前記各待ち行列に結合された手段を更に有する請求項12のルータ。

【請求項14】 前記トラッキング手段は、前記待ち行列に結合されたカウンタデバイスを有し、前記カウンタデバイスは、パケットが前記待ち行列に入力されるたびに増分し、パケットが前記待ち行列から解放されるときは減分する請求項13のルータ。

【請求項15】 前記コントロール手段は、

(a) 前記各待ち行列の現在の待ち行列長さの値 q_i を取得する手段と、

(b) 現在の待ち行列長さの値 q_i と、各待ち行列について割当てられたバッファ占有率 b_i との間の差を計算する手段とを有し、計算された差で最大の値を有する前記待ち行列を選択する請求項 13 のルータ。

【請求項 16】 前記コントロール手段は、

(a) 現在の待ち行列長さの値 q_i が各待ち行列について割当てられたバッファ占有率 b_i を越える一連の待ち行列を計算する手段と、

(b) 前記一連の待ち行列からランダムに待ち行列を選択する手段とを更に有する請求項 13 のルータ。

【請求項 17】 前記選択された待ち行列から解放されるべき前記パケットは、前記待ち行列に内在する最も古いパケットである請求項 12 のルータ。

【請求項 18】 前記スケジューラは異なる所定速度に従って各バッファからパケットを除去し、前記コントロール手段は、

(a) そのそれぞれの待ち行列長さ q_i とその所定のサービス速度から各待ち行列に格納されたパケットにより経験されるキューイング遅延を計算する手段と、

(b) 最も長いキューイング遅延を経験する待ち行列を選択する手段とを更に有する請求項 13 のルータ。

【請求項 19】 1 つ以上のソースからの情報パケットをリンクを介して宛先へ流すことができるフィードバック制御 TCP 接続からトラヒックを搬送する IP ネットワークのためのフロー毎のキューイング装置であり、(a) 複数個の待ち行列に区分けされた所定サイズのバッファと、ここで、各待ち行列は、情報パケットを受信し、一時的に格納するための占有率 b_i が割当てられる、

(b) 所定の速度に従って各バッファからパケットを除去し、また、前記ネットワークを介して前記パケットを伝送するためのスケジューラと、

(c) 前記受信パケットを受信することができ、前記待ち行列が使用可能である場合、前記パケットを前記待ち行列に入力することができる前記バッファ内の待ち行列の使用可能性を決定するためのコントロールデバイスとからなり、前記コントロールデバイスは更に、最長待ち行列ファーストスキームに従って待ち行列を選択し、そして、前記待ち行列が使用可能でない場合、前記受信パケットの入力を可能にするために前記選択された待ち行列からパケットをドロップし、これにより、前記リンクによる公平性を高め、パケットスループットを高めることを特徴とするフロー毎のキューイング装置。

【請求項 20】 1 つ以上のソースからの情報パケットをリンクを介して宛先へ流すことができるフィードバック制御 TCP 接続からトラヒックを搬送する IP ネットワークのためのフロー毎のキューイング方法であり、(a) 複数個の待ち行列に区分けされた所定サイズのバッ

ファを提供するステップと、ここで、各待ち行列は、情報パケットを受信し、一時的に格納するための占有率 b_i が割当てられる、

(b) 所定の速度に従って各バッファからパケットを除去し、また、前記ネットワークを介して前記パケットを伝送するためのスケジューラを提供するステップと、

(c) 前記受信パケットを受信することができ、前記待ち行列が使用可能である場合、前記パケットを前記待ち行列に入力することができる前記バッファ内の待ち行列の使用可能性を決定するためのコントロールデバイスを提供するステップとからなり、前記コントロールデバイスは更に、最長待ち行列ファーストスキームに従って待ち行列を選択し、そして、前記待ち行列が使用可能でない場合、前記受信パケットの入力を可能にするために前記選択された待ち行列からパケットをドロップし、これにより、前記リンクによる公平性を高め、パケットスループットを高めることを特徴とするフロー毎のキューイング方法。

【発明の詳細な説明】

【0001】

【発明の属する技術分野】 本発明は TCP プロトコルを処理する通信ネットワークに関する。更に詳細には、本発明はフィードバック制御ネットワークトラヒック用の接続ごとのキューイング (queue) 及びバッファ管理スキームに関する。

【0002】

【従来の技術】 データパケットの高信頼トランスポートを実行するためのトランスポートプロトコル "TCP" 処理系は周知である。図 1 は、TCP/IP ネットワーク (例えば、インターネット、イントラネット) における従来の技術によるフロー毎の装置の概要図である。

【0003】 図 1 において、ローカル・エリア・ネットワーク (LAN) 20 はルータ 50 を介して広域ネットワーク (WAN) 60 に接続されている。一般的に、ルータ 50 のアーキテクチャは、複数のデータストリーム (例えば、LAN 20 に接続されたソースデータ端末からのもの) が単一のバッファにより受信され、かつ、リンク 55 を介してパケットを送信するサーバにより処理されるようなものである。

【0004】 データパケットがその宛先 (例えば、広域ネットワーク 60 に配置されたデータ端末) に到着すると、パケットの受信機は、受信された各データパケットについて、暗黙的に又は明快に、肯定応答 "ACK" パケットを発生する。ソース端末からのデータ送信は、先行データ送信からの受信肯定応答パケットの速度により制御される。

【0005】 従って、周知のように、ソースからのデータパケット送信の制御は、Richard W. Stevens, "TCP/IP Illustrated", Vols. 1-3, Addison Wesley 1994 に概説されるように、TCP-Reno と TCP-Tah o

eの、適応型スライディングウィンドウスキームの2つの主要変形に従い、受信機からのデータACKパケットのフィードバックにより決定される。

【0006】TCPは長ラウンドトリップ時間を伴う接続に対して固有の不公平性(unfairness)を示すことが確認される。すなわち、短ラウンドトリップ時間接続は、長ラウンドトリップ時間接続よりも大きな帯域幅を得てしまう。なぜなら、短ラウンドトリップ時間接続の場合のTCPウィンドウは、長ラウンドトリップ時間接続のウィンドウよりも速く成長してしまうからである。これは、短ラウンドトリップ時間接続では接続肯定応答(ACK)がより速く受信されるためである。この固有の不公平性を解消する効果的な解決策は未だ発見されていない。

【0007】TCPに固有の別の問題点は、ウィンドウ同期(window synchronization)現象である。例えば、2つの等しいラウンドトリップ時間接続がリンクにあるとすれば、そのリンクは飽和してしまい、バッファに蓄積が始まる。或る時点で、そのリンクに接続されたバッファはオーバフローしてしまい、そして、これらの接続はおおよそ同時にパケット損失を経験する。結果的に、各接続はそのウィンドウを事実上同時に縮小し、ネットワークスループット速度も同時に低下してしまい、その結果、ネットワークリンクは十分に活用されない。

【0008】TCP接続ネットワーク性能を低下させるその他の問題点は、位相効果、ボトルネックリンクによる低リバースパスの物理的な帯域幅非対称性、バーストラヒックの生成、分離性の欠如、より侵略的なトランスポートプロトコル又は悪意のユーザからの保護性の欠如などである。

【0009】ネットワーク性能を改善するために実行された或る解決法は、ランダム・アーリー・ディテクション("RED")スキームである。このスキームは、バッファが飽和される前に、パケットをドロップするように機能する。バッファはモニターされ、平均バッファ占有率が特定の閾値を越えたら、パケットは、例えば、平均バッファ占有率(移動時間平均)の関数である特定の確率に従ってランダムにドロップされ、結果的に、TCP接続はそのウィンドウをランダムにカットする。

【0010】これは、ウィンドウを非同期化するのに特に有効である。特に、FIFO-REDスキームでは、リンクを多数使用する接続は、一層多くのパケットをドロップし易い。従って、結果的に、その対応ウィンドウサイズもカットされる。これは、このような高速度接続のバイアスを低下させる傾向がある。

【0011】米国特許出願第08/858310号明細書に開示された別のスキームは、リバースパスが混雑した時に、公平性とスループットを高める効果をもつACKパケットに関するフロントストラテジーからのドロップを実行する。

【0012】非TCP再ルーブリック(例えば、いわゆる"リーキー・バケツ(leaky-bucket)"条件付タイプのストリーム)の場合、等しい又は異なる重みに従い、各入力接続間で帯域幅が共有されるように実行されるスケジューリングスキームに従って、接続リンクの性能を最大化するための公平待ち行列スケジューリングスキームが実行されている。

【0013】最近の研究努力は、エンド・ツー・エンド遅延のバウンド(bound)及び分離(isolation)を保証する手段として、公平待ち行列スキームを使用することに集中的に向けられている。これに関連する研究は基本的に、非フィードバック制御リーキーバケット管理トラヒックに関するものであり、重み付公平キューイングスキームによる交換機及びルータが発達した。

【0014】例えば、A.K. Parekh and R.G. Gallager, "A Generalized Processor Sharing Approach to Flow Control in the Single Node Case", Proc. of INFOCOM '92, pp. 915-924, May 1992及びD. Stiliadis and A. Varma, "Design and Analysis of Frame-based Queueing: A New Traffic Scheduling Algorithm for Packet-Switched Networks", Proc. of ACM SIGMETRICS '96, p. 104-115, May 1996参照。

【0015】交換機における公平キューイング(FQ)の関心は、様々なバッファ構成(例えば、物理的又は論理的に分離されたバッファ、共用バッファなど)のようなバッファ構成)を処理するためのスケジューラの使用に集中している。

【0016】

【発明が解決しようとする課題】トラヒックを制御し、フィードバック制御TCPネットワークの性能を改善するための公平キューイングスキームの実現が強く望まれている。

【0017】更に、TCP接続に関する相当のスループットを可能にするパケットドロップ機構を実行する公平キューイングスキームの実現が強く望まれている。

【0018】

【課題を解決するための手段】前記課題は、本発明による、フィードバック制御TCPネットワークにおける公平キューイングスキームと共に併用されるフロー/接続毎の共用バッファ管理スキームにより解決される。

【0019】本発明によれば、長ラウンドトリップ時間を伴う接続に対するTCPの固有の不公平性を軽減する、異なるTCPバージョンを用いる接続がボトルネックリンクを共有する場合、分離性を与える、一層侵略的なトラヒックソース、不品行なユーザ又はリバースパス輻輳の場合のその他のTCP接続からの保護性を与える、両方向トラヒックの存在下におけるACK圧縮作用を軽減する、ACK損失(トラヒックのバーストを起す)を経験するユーザに顕著な悪影響を及ぼすその他の接続から保護する、及び全体的なリンク利用率

を低下させることなく、「欲張り」接続を伴うボトルネックを共用する双方向接続に対する低待ち時間を与える、などのようなTCPの種々の目標を達成することができる。

【0020】更に詳細には、共用バッファアーキテクチャは、各接続のための帯域幅予約保証 r_i により実行される。速度 r_i のそのバッファを完全に使用する第1の接続と、十分に利用されない（無駄使いされる）、速度 r_2 の第2の接続があり、第1の接続が更に一層の帯域幅を必要とする場合、共用バッファスキームでは、第1の接続からバッファリング（帯域幅）を借りることもできる。

【0021】第2の接続がそのバッファスペースを満足しなければならぬ場合、別の利用バッファからのデータを押し出し、着信データパケットのための場所を空けなければならない。

【0022】接続毎の待ち行列（キュー）管理スキームは、共用バッファアーキテクチャにおけるバケットドロップングメカニズム（例えば、最長待ち行列ファースト（“LQF”））をサポートし、FIFO-REDバッファ管理スキームよりも優れたTCP性能を生じる。本明細書では、公平性は、各リンク容量に対する個々のスループットの割合の平均値に対する標準偏差の比により計算する。

【0023】

【発明の実施形態】図2は、様々なソース s_1, \dots, s_n から発信するバケットトラヒックを処理するTCPネットワーク接続のルータ65のための接続毎の待ち行列アーキテクチャを示すブロック図である。

【0024】ネットワーク接続要素は、汎用の共用バッファメモリBを有する。このメモリBは、速度Cでリンク55上のデータパケットを処理するスケジューラ75と単一（ボトルネック）のリンク55との接続のために、“i”個の複数の待ち行列30a, ..., iに分割されている。各バッファ接続iは、接続iの保証バッファサイズである各バッファ割当て b_i を有する。

“フロー毎のキューイング”又は“接続毎のキューイング”として知られているこのアーキテクチャでは、待ち行列に入れられるべき到着パケットのきめ細かい動的分類が必要である。

【0025】“ソフト状態”アプローチは、潜在的に非常に多数のおそらく活性な接続（この場合、大多数の接続は実際には活性ではないであろう）へ適く接続状態を維持する。ネットワークノード内の関連状態はその他のガーベジ・コレクション（ごみ集め）の手段によりタイムアウト（時間切れ）を起こさせずか、矯正させるか又は遅延させるために待機される。従って、スケーラビリティ（scalability）及び汎用性は、接続毎のサーバーのための基本的な要件である。必要な全ての操作は、0

(1) 複雑さにより実行可能であり、リソース（バッ

ファ又は帯域幅）は所定の接続には静的に割当てられない。

【0026】操作中、スケジューラ75は、各待ち行列について同等であるか又は所定の重みに応じた、 r_i に等しい速度で各個別待ち行列iを処理する。速度 r_i で広い帯域幅を使用する特定の待ち行列（接続）（例えば、待ち行列30a）は、その他の待ち行列よりも長い待ち行列を有する傾向がある。全ての待ち行列接続が、それらの各帯域幅割当てを完全に利用している場合、待ち行列30aはおそらくバッファオーバーフローを経験するであろう。

【0027】或る割当てメモリ（例えば、30c）が完全に利用されていない場合、本発明の公平キューイング（FQ）スキームは、待ち行列30bに到着するデータパケットのために、必要に応じて、十分に利用されていない待ち行列からバッファスペースを利用するか又は借りることができる。従って、バッファ待ち行列30aの予約割当て b_i を超えることができる。

【0028】第1の高速バッファ待ち行列30aのためのパケットを受信する第2のバッファが満杯になると、別の十分に利用されていないバッファ（例えば、待ち行列30i）は、待ち行列30aに宛てられる新たなデータパケットにバッファスペースを貸す。

【0029】同時に2つ以上の待ち行列がバッファオーバーフローを経験することもでき、そのため、十分に利用されていないバッファスペースから借りることもできる。従って、2つ以上の待ち行列がその予約割当て b_i を超えることもできる。本発明は、TCP接続ネットワークのフォワード及びリバース接続の両方に適用可能であり、また、データ及びACKトラヒックフローの両方に同等に適用可能である。

【0030】接続iが (b_i+1) 個以上のバッファを必要とする場合、総占有率がB未満であるとすれば、使用可能なプールからスペースが割当てられる。

【0031】図3(a)は、本発明の接続毎の待ち行列スキームを実行するための一般的な流れ図である。ステップ201に示されるように、接続iについて宛てられたパケットが到着すると、最初のステップ203で、バッファBの現に残っている使用可能なスペースを有するカウンタをチェックし、新たに到着したパケットを確実に収容するための十分なメモリスペースがバッファ内に残存しているか否か決定する。

【0032】新たに到着したパケットを確実に収容するための十分なメモリスペースがバッファ内に残存している場合、ステップ208で処理が継続され、到着パケットが属する接続iを識別し、そして、ステップ211で、このパケットを、この接続に対応する待ち行列に格納する。このスキームに潜在的に含まれることは、到着パケットが適正に分類され、そして、待ち行列のうちの一つに割当てられることである。

【0033】ステップ203において、新たに到着したパケットを確実に収容するための十分なメモリスペースがバッファ内に残存していないと決定される場合（例えば、現在の占有率 q_i が b_j 未満である待ち行列30jがバッファを必要とする場合）、ステップ225において追出しスキームが実行され、到着パケットのために場所を空ける。特に、呼出されたTCP/IPプロトコルの実行に応じて、追出しが行われる待ち行列を選択するための2つの方法が使用できる。

【0034】特に、図4(b)に示されるように、第1の実施態様における追出しメカニズムは、LQFスキームであり、別の待ち行列から予約された最大量のメモリを借りる待ち行列、すなわち、 $(q_i - b_i)$ が全ての接続を通して最大であるような接続を選択する。例えば、ステップ250に示されるように、バッファ内の各待ち行列の現在のバッファ割当てに関する決定が行われる。

【0035】次いで、ステップ260で、各待ち行列の現在の待ち行列長さ q_i を取得し、ステップ270で、各待ち行列の $(q_i - b_i)$ 差に関する計算を行う。最後に、ステップ275において、最大の $(q_i - b_i)$ 差を有する待ち行列を選択する。従って、その予約割当て b_i からの最大偏差は最長待ち行列であり、このため、図3(a)のステップ226で示されるように、到着パケットと先入れ相関関係を有する待ち行列からパケットがドロップされる。

【0036】当業者なら前記の最長待ち行列を最初に追出すスキーム(LQFスキーム)を実行するためのその他のアルゴリズムを案出することもできるので、本発明は図4(b)に示された方法論に限定されない。例えば、最長遅延ファースト("LDF")ドロッピングメカニズムは、割当て処理速度 r_i が全て等しい場合、LQFスキームと同等に実行できる。なぜなら、待ち行列が同じ速度で処理される場合、遅延は各接続について同一だからである。同様に、処理速度が均でない場合、待ち行列長さが同一であっても、遅延は異なる。従って、LQFスキームはLDFの特例である。

【0037】第2の最長待ち行列よりかなり長い待ち行列を有する多数の接続を有するシステムで実行する場合、最長待ち行列ファースト(LQF)は過度のバースト損失をもたらすことがある。すなわち、2個以上の非常に接近したパケットは、その割当てを越える待ち行列から連続的にドロップされる。例えば、TCP-Reno型処理系はバースト損失の存在下で不良動作を行うことが知られているので、TCP-Renoアーキテクチャを実行する接続の性能は悪影響を受ける。従って、前記LQFにおけるバースト損失量を低下するために、スキームは、それぞれの割当てを越えるバッファからランダムに抽出するランダムジェネレータを使用するように変更されている。

【0038】特に、第2の実施態様では、各バックログ接続(backlogged connection:予備接続) i は、 $b_i = B/n$ (ここで、 n はバックログ接続の個数である) で示される各目バッファ割当てを有する。図5

(c)に示されるように、ステップ255で、各待ち行列のメモリ割当て b_i を取得する。次いで、ステップ265で、バックログ接続を例えば、2個のサブセットにグループ分けする。一方は、 b_i よりも大きな占有率 q_i を有するグループであり、他方は、 $q_i \leq b_i$ の関係を有するグループである。

【0039】割当てより上(すなわち、 $q_i > b_i$)の待ち行列群から、ステップ285で示されるように、1つの待ち行列がランダムに選択され、そして、ステップ226で示されるように、1つのパケットが先頭からドロップされる。

【0040】異なる接続に関するバッファ占有率が平坦化し、そして、システムに過負荷をかけられる接続から最適に保護するための試みにおいて、L. Geraciadis, I. Ci don, R. Guerin, and A. Khamis, "Optimal Buffer Sharing," IEEE J. Select. Areas Commun., vol. 13, pp. 1229-1240, Sept. 1995に記載されるような開ループトラフィックについて実行されるスキームと同様な方法で、特定の追出し(pushout)スキームは最長待ち行列の先頭からパケットをドロップする。

【0041】図6に示されるように、ハードウェアの観点から、カウンタ90a, ..., 90iは、それぞれの付随待ち行列の占有率 q_i のトラックを保持するようプログラムされたコントロールプロセッサ92と共に、各待ち行列30a, ..., 30iに関連するように図示されている。従って、図4(b)及び図5(c)のステップ260及び265において、現在の待ち行列長さは、例えば、ポーリングにより、又は、図6のテーブル99のようなレジスタテーブル(ここで、レジスタは最長待ち行列長さを示す)から位置指定することにより取得される。

【0042】別のカウンタ95は、総バッファ占有率 B のトラックを保持するために図示されている。従って、パケットが到着すると、プロセッサ92は、カウンタ95で使用可能な総メモリを決定するために、ステップ203(図3(a)参照)におけるチェックを行う。

【0043】最長待ち行列を決定するために、プロセッサ92はソート済み構造を実行することもできる。例えば、各キューに入れる、キューから外す又はドロップ操作の時点で、待ち行列の対応カウンタがそれに応じて増分又は減分された後、最長待ち行列構造を常に見ておくために、その待ち行列占有率値 q_i を現在の最長待ち行列占有率値と比較するように、ソート済み構造を実行する。

【0044】これらフロー毎のスキームの変法は、指数関数的に増大するサイズを有する一連のビンを使用する

ことにより効率的に実行することができる。待ち行列が変更される場合(例えば、キューに入れる、キューから外す又はドロップする場合)は常に、システムは最高占有ビンのトラックを保持しながら、適当なビン(例えば、現在のビンと同一のビン、一つ上のビン又は一つ下のビン)に移動させる。

【0045】バッファが満杯の場合、最高占有ビンからの待ち行列が選択され、その最初のパケットがドロップされる。バッファがバイト単位で測定される場合、この操作は、可変パケットサイズのため、新たに到着したパケットを収容するために十分なスペースが解放されるまで、繰り返さなければならない。真のLQDを行うために、最高占有ビンにおける待ち行列をソート済みリストとして維持するが、又はいつも最長待ち行列をサーチしなければならない。

【0046】FIFO-RED及びFQ-REDスキームと比較される本発明のバッファ管理スキームの改善された性能のシミュレーションについて説明する。図7は、このシミュレーションシステム119を示す概要図である。図7において、ソース101a、・・・、101fはルータ120への高速アクセスバスを有する。ルータ120は唯一のボトルネックであり、本発明の接続/フロー毎のバッファ管理スキームを実行するものである。

【0047】比較は、TCP-Tahoe及びTCP-Renoソース、パースト及び欲張りソース(greedy source)、リバースバス転換を有する片方向及び両方向トラフィックソース、及び広範囲に異なるラウンドトリップ時間の混合を用いて行われる。アクセスバス遅延は、様々なラウンドトリップ時間のモデルを作るために広範囲にわたって設定され、宛先はパケット毎にACKと仮定した。

【0048】片方向トラフィックの場合、ACKは非輻輳パスを介して送信される。両方向トラフィックの場合、リターンパスはルータを経由し、そのため、キューイング(待合わせ)遅延が生じることがある。特に、ルータがFIFOスケジューリングを使用する場合、ACK及びデータパケットは待ち行列内で混合される。公平キューイングの場合、ACKは別のフローとして処理される。非対称トラフィックの場合、リターンリンクの帯域幅は、宛先からルータまで低下される。そのため、相当なリバースバス転換とACK損失が生じる。

【0049】シミュレーションシステム119におけるTCP-Tahoe及びTCP-Renoの処理系は、4.3-Tahoe BSD及び4.3-Reno BSDのTCPフローと輻輳制御挙動とをそれぞれモデル化する。REDモデルは有向パケットであり、最小閾値としてバッファサイズの25%を使用し、最大閾値としてバッファサイズの75%を使用する。待ち行列重みは0.002である。

【0050】図8(a)及び(b)は、従来技術によるFIFO-RED及びLQF法と比較されるような、利用率及び公平性のそれぞれの観点から、擬似TCPネットワークで公平待ち行列LQD及びRNDドロップ法を実行した場合の、改善された性能を示す特性図である。

【0051】特に、図8(a)及び(b)は、図7の擬似ネットワークにおける10Mbps/100kbps容量(TCP-Tahoe及びTCP-Renoの20ms-160msラウンドトリップ時間)を有する非対称ボトルネックリンクによる20個のTCP接続のためのバッファサイズ(で表したリンク速度)の関数として、スループット(図8(a))及び公平性係数(図8(b))により評価される改善されたTCP性能を示す。

【0052】図示されているように、ライン137a及び138aでそれぞれ示されるFQ(公平キューイング)-RED及びFIFO-REDの両戦略は、ライン139及び140で示されるFQ-LQF及びFQ-RNDドロップ法よりも非常に低いスループットを有する。なぜなら、再送パケットに対応するACKが、擬似非対称値が3(これは帯域幅非対称ではない)の場合の時間が66%も喪失されるからである。これは、TCPサイクルの少なくとも66%でタイムアウト(時間切れ)を生じ、スループットを著しく低下させる。

【0053】フォワードパスにおける多数の損失及びフォワードパスにおける再送パケットの損失のために、その他のタイムアウトも起こる。一方、リバースバスにおける先頭からのドロップはこれらのタイムアウトを殆ど完全に除去する。タイムアウトは不経済なので、両方のREDスキームとも、FIFO-LQDを含むその他のスキームよりも著しく低いスループットを有する。

【0054】更に、図8(b)に示されるように、FQ-RND及びFQ-LQDワークの両方とも非常に良好である。なぜなら、これらは、フロー毎の待ち行列の効果と、先頭からのドロップのタイムアウト除去効果とを併有するからである。FQ-LQDは、再送パケットのドロップに対する組み込みバイアスを有するという別の効果も有する。これは、ソースが最初の重複ACKの受信による損失を検出したときに、パケットの送信を中止するからである。

【0055】第3の重複ACKが受信された後にだけ、再送パケットが送信される。介入間隔中、ソースがTCPにより強制的に不活動化される場合、フローに対応する待ち行列は、少なくともその最小保証速度で排出される。従って、再送パケットが到着する場合、最長待ち行列であることは殆どない。従って、再送パケットのドロップに対する固有バイアスである。このバイアスは非対称ネットワークに限定されないが、低速リバースチャネル膨張のために、非対称要因、第1の重複ACK受信と第3の重複ACK受信との間の間隔により、バイア

スは非対称ネットワークにおいて高められる。

【0056】再送パケットの損失は不経済なタイムアウトを起こすので、このパイアスは、図8(b)においてライン141で示されるように、FQ-LQDの性能を改善する。ライン142で示されるFQ-RNDも同様このパイアスを有するが、著しく低い。その理由は、FQ-LQDに対するものと若干類似している。

【0057】第1の重複ACK受信と第3の重複ACK受信との間の間隔中に、フローの待ち行列は、(ソースは不活動化されているので)少なくともその保証速度と等しい速度で流出し、待ち行列占有率はこのフローの予約パラメータ以下にまで落ちるであろう。この場合、再送パケットが到着すると、集合バッファが満杯であったとしても、再送パケットは失われない。これらの効果より、FQ-LQD及びFQ-RNDは最高の性能を有する。

【0058】

【発明の効果】以上説明したように、本発明によれば、

長ラウンドリップ時間を伴う接続に対するTCPの固有の不公平性を軽減する、異なるTCPバージョンを用いる接続がボトルネックリンクを共用する場合、分離性を与える、一層侵略的なトラフィックソース、不品行なユーザ又はリバースパス輻輳の場合のその他のTCP接続からの保護性を与える、両方向トラフィックの存在下におけるACK圧縮作用を軽減する、ACK損失(トラフィックのバーストを起こす)を経験するユーザに顕著な悪影響を及ぼすその他の接続から保護する、及び全体的なリンク利用率を低下させることなく、“欲張り”接続を伴うボトルネックを共用する双方向接続に対する低待ち時間を与える、などのようなTCPの種々の

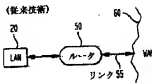
目標を達成することができる。

【図面の簡単な説明】

【図1】従来技術によるTCPネットワーク接続を示すブロック図である。

【図2】多数のTCP接続のための共用バッファアーキ

【図1】



テクチャのブロック図である。

【図3】(a)はバッファ割当て及びパケットドロップングスキームを行うために履行される方法を示す流れ図であり、図3、図4、図5で1つの流れ図を形成する。

【図4】(b)はLQFパケットドロップ法を示す流れ図である。

【図5】(c)はRNDパケットドロップ法を示す流れ図である。

【図6】各バッファ待ち行列に関連するハードウェアを示すブロック図である。

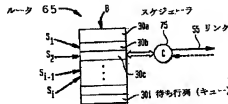
【図7】公平キューイングを処理するルータを有するTCP/IPネットワークのシミュレーションを示すブロック図である。

【図8】(a)は図7の擬似ネットワークにおけるボトルネックリンクによる20個のTCP接続のためのバッファサイズ(で表したリンク速度)の関数として、スループットにより評価される改善されたTCP性能を示す特性図である。(b)は図7の擬似ネットワークにおけるボトルネックリンクによる20個のTCP接続のためのバッファサイズの関数として、公平性係数(全リンク容量に対する個々のスループットの割合の平均値に対する標準偏差の比)により評価される改善されたTCP性能を示す特性図である。

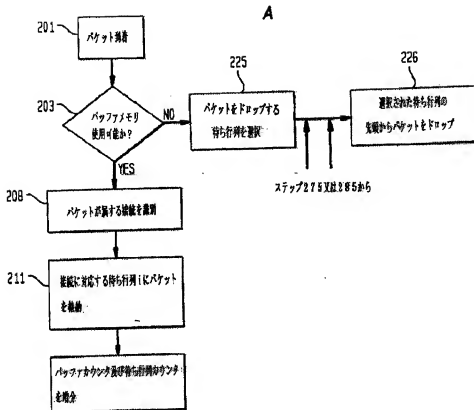
【符号の説明】

- 30 待ち行列 (キュー)
- 55 ルータ
- 65 ルータ
- 75 スケジューラ
- 92 コントロールプロセス
- 95 カウンタ
- 99 レジスタテーブル
- 101 ソース
- 119 シミュレーションシステム
- 120 ルータ

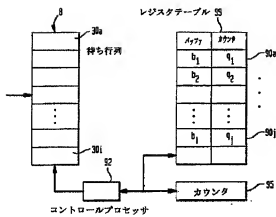
【図2】



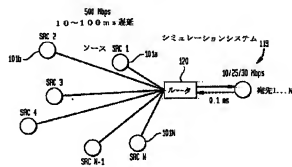
【図3】



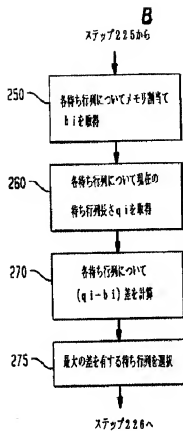
【図6】



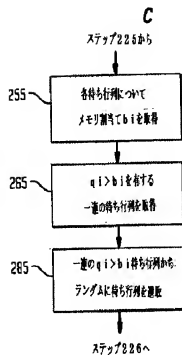
【図7】



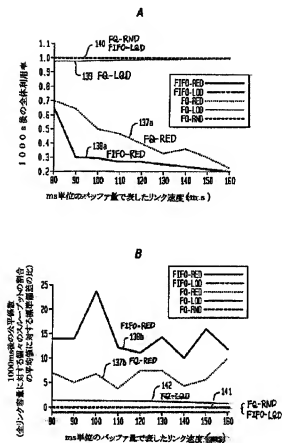
【図4】



【図5】



【図8】



フロントページの続き

(71) 出願人 596077259

600 Mountain Avenue,
Murray Hill, New Je
rsey 07974-0636 U. S. A.

(72) 発明者 ティー、ヴィー、ラクシュマン
アメリカ合衆国, 07724 ニュージャージ
ー, イートンタウン, ヴィクトリア ドラ
イブ 118

(72) 発明者 デイミトリオス スティリアディス
アメリカ合衆国, 07030 ニュージャージ
ー, ホボケン, パーク アヴェニュー
106

(72) 発明者 バーナード シュター
アメリカ合衆国, 07747 ニュージャージ
ー, アバーディーン, アイダホ レイン
12